

SUPPLEMENTARY MATERIAL

This material supplements the paper “Extracting Top- k Frequent and Diversified Patterns in Knowledge Graphs”. It includes pseudocodes, proofs, complexity analysis, and complementary experimental results for kCP and kCPD.

APPENDIX A

PSEUDOCODES

We first present the pseudocodes of solutions for kCP and kCPD in the following.

Algorithm 1, denoted as kCP-B, summarizes the procedures of the baseline solution to kCP problem (see in Section 3.1). It invokes two sub-routines, i.e., pattern extension (Line 9) and MNI evaluation (Line 6), which correspond to the Algorithm 2 and Algorithm 3, respectively.

Algorithm 1 kCP-B(KG, P , k)

```

1: Initialize top- $k$  frequent pattern result min-heap  $\mathcal{H} \leftarrow \emptyset$ 
2: Initialize max-heap candidate pattern set  $\mathcal{CS} \leftarrow (P, 0)$ 
3: while  $\mathcal{CS}$  is not empty do
4:    $P' \leftarrow \mathcal{CS}.\text{top}(); \mathcal{CS}.\text{delete\_max}()$ 
5:   if  $P'.\text{frq} > \mathcal{H}.\text{top}().\text{frq}$  then
6:      $P'.\text{frq} \leftarrow \text{EvaluateMNI}(P', \text{KG})$   $\triangleright$  Algorithm 3
7:     if  $P'.\text{frq} > \mathcal{H}.\text{top}().\text{frq}$  or  $\mathcal{H}.\text{Size} \leq k$  then
8:       Update  $\mathcal{H}$  by  $(P', P'.\text{frq})$ 
9:        $\text{PatExt}(P', \text{KG}, \mathcal{CS})$   $\triangleright$  Algorithm 2
10: return  $\mathcal{H}$ 

```

Algorithm 2 PatExt(P , KG, \mathcal{CS})

```

1: for each  $(e, v) \in E \times V$  do
2:   if  $e$  can be used to extend  $P$  then  $\triangleright$  gSpan
3:     Let  $P'$  be the extension of  $P$  with  $e$ 
4:      $P'.\text{frq} \leftarrow P.\text{frq}$ 
5:      $\mathcal{CS} \leftarrow \mathcal{CS} \cup \{(P', P'.\text{frq})\}$ 

```

Algorithm 3 EvaluateMNI(P , KG)

```

1:  $\text{count} \leftarrow +\infty$ 
2: for each node  $v \in V_P$  in pattern  $P$  do
3:    $\text{tmp} \leftarrow 0$ 
4:   for each candidate instance  $i_v$  of  $v$  do
5:     if  $i_v$  has been marked then
6:        $\text{tmp}++$ 
7:     else if  $\exists$  instance  $S$  of  $P$  and  $S$  includes  $i_v$  then
8:        $\mathcal{X}(P) \leftarrow \text{mark } S\text{'s nodes to } V_P\text{'s domains}$ 
9:        $\text{tmp}++$ 
10:  if  $\text{tmp} < \text{count}$  then  $\text{count} \leftarrow \text{tmp}$ 
11: return  $\text{count}$ 

```

Algorithm 4 presents the pseudocode of pattern enumeration via meta-index in Section 3.2.1.

Algorithm 4 PatExtMeta(P , M , \mathcal{CS})

```

1: for each node  $u \in V_P$  do
2:   for each  $(e, v) \in M(u).\text{Out} \cup M(u).\text{In}$  do
3:     if isValid( $P$ ,  $(e, v)$ ) then  $\triangleright$  node degree checking
4:       if  $P$  can be extended by  $(e, v)$  then  $\triangleright$  gSpan
5:         Let  $P'$  be the extension of  $P$  with  $(e, v)$ 
6:          $P'.\text{ub} \leftarrow P.\text{frq}$ 
7:          $\mathcal{CS} \leftarrow \mathcal{CS} \cup \{(P', P'.\text{ub})\}$ 

```

The algorithm below (i.e., Algorithm 5) shows the pseudocode of join⁺ MNI computation, which is a core sub-routine of FastPat framework (see in Section 3.2.4).

Algorithm 5 JoinMNI⁺(P , PT, θ)

```

1: Initialize join result table  $R \leftarrow \emptyset$ 
2: Serialize pattern  $P$  into a set of path patterns  $\mathcal{PS}$ 
3: for each path pattern  $\mathcal{P} \in \mathcal{PS}$  do
4:    $\text{PT} \leftarrow \text{prune tuples from PT by table join on } \mathcal{P}$ 
5: Apply tuple reduction on each edge table of PT
6:  $\text{MNI}_{lb}(P) \leftarrow \text{compute the MNI lower bound of } P$ 
7: while  $\mathcal{PS}$  is not empty do
8:    $\mathcal{PS} \leftarrow \mathcal{PS} - \{\mathcal{P} \leftarrow \text{argmin}_{\mathcal{P} \in \mathcal{PS}} \text{MNI}_{lb}(P)\}$ 
9:    $\mathcal{T}_{\mathcal{P}} \leftarrow \text{edge tables of path pattern } \mathcal{P} \text{ from PT}$ 
10:  while  $\mathcal{T}_{\mathcal{P}}$  is not empty do
11:     $\mathcal{T}_{\mathcal{P}} \leftarrow \mathcal{T}_{\mathcal{P}} - \{\mathcal{T}(A, B) \leftarrow \text{argmin}_{\mathcal{T} \in \mathcal{T}_{\mathcal{P}}} \text{MNI}(\mathcal{T})\}$ 
12:     $R \leftarrow \text{join } \mathcal{T}(A, B) \text{ with } \forall \mathcal{T}_i \in \text{PT} - \{\mathcal{T}(A, B)\}$ 
13:     $P.\text{ub} \leftarrow \text{the MNI of verified } \mathcal{T}_{\mathcal{P}} \text{ in } R$ 
14:    if  $P.\text{ub} < \theta$  then  $\triangleright$  Pruning by MNI upper bound
15:       $P.\text{frq} \leftarrow P.\text{ub}$  and return
16:  $P.\text{frq} \leftarrow \text{MNI}(R)$ 

```

Algorithm 6 is the pseudocode of the baseline solution (i.e., FPD-B) for kCPD problem (see in Section 4.1).

Algorithm 6 FPD-B(KG, P , k , λ)

```

1: Initialize max-heap candidate pattern set  $\mathcal{CS} \leftarrow \{(P, 0)\}$ 
2: Initialize top- $k$  pattern result set  $\mathcal{R} \leftarrow \emptyset$ 
3: while  $|\mathcal{R}| < k$  and  $\mathcal{CS}$  is not empty do
4:    $\mathcal{CS}.\text{delete}(P' \leftarrow \mathcal{CS}.\text{top}())$ 
5:    $\mathcal{E} \leftarrow \text{PatExt}(P', \text{KG}, \emptyset)$   $\triangleright$  Algorithm 2
6:   for each candidate pattern  $d \in \mathcal{E}$  do
7:      $d.\text{frq} \leftarrow \text{EvaluateMNI}(d, \text{KG})$   $\triangleright$  Algorithm 3
8:      $\mathcal{CS} \leftarrow \mathcal{CS} \cup \{(d, d.\text{frq})\}$ 
9:    $P'' \leftarrow \mathcal{CS}.\text{top}()$ 
10:  if  $|\mathcal{R}| = 0$  or  $\forall p \in \mathcal{R}, \text{Div}(p, P'') \geq \lambda$  then
11:     $\mathcal{R} \leftarrow \mathcal{R} \cup \{P''\};$ 
12: return  $\mathcal{R}$ 

```

APPENDIX B

PROOFS

We next give the proofs of lemmas in the paper.

Proof of Lemma 1. For each node $v \in V_P$, if the node v 's instance u is not in $V_{\text{Dom}}(v)$, it also cannot be in $F(v)$, the

image set of v holds $F(v) \subseteq \text{VDom}(v)$. We have:

$$\begin{aligned} & \forall v \in V_P, F(v) \subseteq \text{VDom}(v) \\ \Leftrightarrow & \forall v \in V_P, |F(v)| \leq |\text{VDom}(v)| \\ \Leftrightarrow & \min_{v \in V_P} |F(v)| \leq \min_{v \in V_P} |\text{VDom}(v)| \end{aligned}$$

Hence, it holds $\text{MNI}(P) \leq \text{MNI}_{ub}(P)$.

Proof of Lemma 2. We construct table $T'(A, B)$ as follows: (i) group all tuples in $T(A, B)$ according to distinct values in $\text{Dom}(T.A)$ and form $\text{MNI}(T.A)$ groups, then randomly sample one tuple from each group and insert it into $T'(A, B)$; (ii) group all tuples in $T(A, B)$ by the distinct values in $\text{Dom}(T.B)$ and form $\text{MNI}(T.B)$ groups, randomly sample one tuple per group and insert it into $T'(A, B)$ until $\text{MNI}(T'.B) = \text{MNI}(T.A)$. We have $\text{MNI}(T') = \text{MNI}(T) = \text{MNI}(T.A)$ as step (i) and (ii) ensure $|\text{Dom}(T'.A)| = \text{MNI}(T.A)$ and $|\text{Dom}(T'.B)| = \text{MNI}(T.A)$, respectively. The total number of tuples in $T'(A, B)$ is bounded by $2 * \text{MNI}(T.A)$ as both step (i) and step (ii) insert at most $\text{MNI}(T.A)$ tuples.

Proof of Lemma 3. It holds as every tuple that appears in the join result of the reduced tables (i.e., R) is included in the exact join result.

Proof of Lemma 4. First, given a graph $G(V, E)$ and an integer n . The decision version of Maximum Independent Set (MIS) problem asks whether there exists a maximum independent set $I = \{v_1, v_2, \dots, v_m\}$ of node set V , such that $|I| \geq n$. We denote MIS problem instance as $\langle G(V, E), n \rangle$, this problem is shown to be NP-hard [1].

Second, consider a knowledge graph KG , a core pattern P , a diversity threshold λ and an integer k . Let V_P includes all the patterns extended from core pattern P , i.e., $V_P = \{P_1, P_2, \dots, P_m\}$. For the sake of NP reduction, we set $\text{MNI}(P_i) = 1$ for each candidate pattern P_i in V_P . Then we transform the kCPD problem by creating an edge set E_P . For every two candidate patterns P_i and P_j in V_P , there is an edge e_{ij} in E_P iff $\text{Div}(P_i, P_j) < \lambda$. In other words, we connect these similar candidate patterns by an edge, thus, the diversified patterns are not connected. Until now, we obtain a graph $G_P(V_P, E_P)$. Thus, the decision version of our simplified and transformed kCPD problem is whether there exists a set S which includes k independent nodes from $G_P(V_P, E_P)$. We denote the problem instance as $\langle G_P(V_P, E_P), k \rangle$.

Third, reducing the MIS problem instance $\langle G(V, E), n \rangle$ to our simplified and transformed kCPD problem instance $\langle G_P(V_P, E_P), k \rangle$ is straightforward, i.e., setting $k = n$, and mapping every node in V to V_P and every edge in E to E_P . It costs polynomial time, i.e., $O(|V| + |E|)$. If there exists an independent set I of size k in MIS problem instance $\langle G(V, E), n \rangle$, then there exists a size- k set S for our simplified and transformed kCPD problem instance $\langle G_P(V_P, E_P), k \rangle$. The reason is that every two nodes v_i, v_j in I correspond to two diversified patterns P_i, P_j in S with diversity score $\text{Div}(P_i, P_j) \geq \lambda$. Conversely, it also holds.

Proof of Lemma 5. We analyze the approximate ratio of the greedy algorithm FPD-B (i.e., Algorithm 6) as follows. First, we denote P^* as the most frequent pattern extended from the core pattern in KG . Let \mathcal{R} and \mathcal{R}' denote the result

output by algorithm FPD-B and an exact algorithm, respectively. Next, we require k patterns returned by the exact algorithm, i.e., $|\mathcal{R}'| = k$. Then, there would be two possible cases for the greedy algorithm FPD-B: (i) It identifies pattern sets with the largest cardinality while being smaller than k , i.e., $|\mathcal{R}| < k$. Since \mathcal{R} always includes pattern P^* , the approximate ratio:

$$\alpha = \frac{\sum_{P_j \in \mathcal{R}'} \text{MNI}(P_j)}{\sum_{P_i \in \mathcal{R}} \text{MNI}(P_i)} \leq \frac{k \cdot \text{MNI}(P^*)}{1 \cdot \text{MNI}(P^*)} = k.$$

(ii) It identifies the set with k patterns that has the largest total frequency among all pattern sets, i.e., $|\mathcal{R}| = |\mathcal{R}'| = k$, we have the approximate ratio of algorithm FPD-B is

$$\begin{aligned} \alpha &= \frac{\sum_{P_j \in \mathcal{R}'} \text{MNI}(P_j)}{\sum_{P_i \in \mathcal{R}} \text{MNI}(P_i)} \leq \frac{|\mathcal{R}'| \cdot \max_{P_j \in \mathcal{R}'} \{\text{MNI}(P_j)\}}{|\mathcal{R}| \cdot \min_{P_i \in \mathcal{R}} \{\text{MNI}(P_i)\}} \\ &= \frac{|\mathcal{R}'|}{|\mathcal{R}|} \cdot \frac{\text{MNI}(P^*)}{\min_{P_i \in \mathcal{R}} \{\text{MNI}(P_i)\}} = \frac{\text{MNI}(P^*)}{\min_{P_i \in \mathcal{R}} \{\text{MNI}(P_i)\}}. \end{aligned}$$

Note that both approximation ratios are not constant. The first case shows the approximate ratio for the greedy algorithm can be k in the worst case. However, k is usually not large and the worst case seldomly happened in practice.

Proof of Lemma 6. By noting that the relations $\mathcal{M}_P \subseteq \mathcal{C}_P \subseteq \mathcal{D}_P$, we have

$$\begin{aligned} \text{Div}(P_1, P_2) &= 1 - \frac{|\mathcal{C}_{P_1} \cap \mathcal{C}_{P_2}|}{|\mathcal{C}_{P_1} \cup \mathcal{C}_{P_2}|} \geq 1 - \frac{|\mathcal{C}_{P_1} \cap \mathcal{D}_{P_2}|}{|\mathcal{C}_{P_1} \cup \mathcal{C}_{P_2}|} \\ &\geq 1 - \frac{|\mathcal{C}_{P_1} \cap \mathcal{D}_{P_2}|}{|\mathcal{C}_{P_1} \cup \mathcal{M}_{P_2}|} = \text{Div}_{lb}(P_1, P_2). \end{aligned}$$

Hence, it holds that $\text{Div}(P_1, P_2) \geq \text{Div}_{lb}(P_1, P_2)$.

APPENDIX C COMPLEXITY ANALYSIS

In this part we analyze the time and space complexities of the proposed algorithms in kCP and kCPD.

C.1 Complexity Analysis of kCP Algorithms

Time Complexity of Algorithm kCP-B. The time cost is dominated by two core sub-routines: pattern enumeration (Algorithm 2) and MNI computation (Algorithm 3). Given a knowledge graph $KG(V, E)$, a core pattern $P(V_P, E_P)$, the number of all possible candidate pattern is $2^{|E|}$, as analyzed in [2]. For each pattern, the time complexity to compute its MNI value is $O((|V_P| + k) \cdot |V|^{|V_P| + k - 1})$ according to the time cost analysis in GRAMI [3]. Thus, the total time complexity of kCP-B is $O(2^{|E|} \cdot (|V_P| + k) \cdot |V|^{|V_P| + k - 1})$.

Time Complexity of Algorithm kCP-A. It consists of three parts: meta-index and edge table construction, pattern enumeration, upper bound and join-based MNI computation. We analyze the time cost of each part as following. First, the time complexity to construct meta-index and edge table is $O(|V| + |E|)$, as it only scans all nodes and edges once. Second, the number of enumerated possible candidate patterns is $(|V_P| + k) \cdot d_{max}^k$ according to Algorithm 4. In particular, the largest extension pattern of core pattern P has $|V_P| + k$ nodes with the anti-monotonicity property of MNI. For each pattern P' , the number of generated extensions will less

than $(|V_P| + k)d_{max}$. Third, the time complexity to compute the MNI upper bound is $O(|V| + |E|)$ as it only need check all nodes and edges in KG. The time complexity of pattern serialization (Algorithm 5) is $O((|V_P| + k) + (|E_P| + k))$ as it only incurs DFS search once. The time complexity join operators in JoinMNI⁺ (Algorithm 5) is $O(|E|^{(|E_P|+k)})$ as we conduct table join for every edge in the candidate pattern. In summary, the total time complexity of kCP-A is $O(|V| + |E| + (|V_P| + k)^k d_{max}^k \cdot (|V| + |E| + (|V_P| + k) + (|E_P| + k) + |E|^{(|E_P|+k)}))$. It is $O((|V_P| + k)^k d_{max}^k \cdot |E|^{(|E_P|+k)})$ after simplification.

Space Complexity of Algorithm kCP-B. The space cost of kCP-B is contributed by the MNI computation, which is adapted from GRAMI, we refer the interested reader for the reported space cost in [3].

Space Complexity of Algorithm kCP-A. It is dominated by intermediated result table size of join-based MNI computation. In particular, for each candidate pattern, join-based MNI computation joins at most $(|E_P| + k)$ edge tables with the anti-monotonicity property of MNI. For each edge table, the maximum table size is $|E|$, thus, the total space cost is $O(|E|^{(|E_P|+k)})$.

C.2 Complexity Analysis of kCPD Algorithms

Time Complexity of Algorithm FPD-B. The time cost is dominated by three core sub-routines: pattern enumeration (Algorithm 2), MNI evaluation (Algorithm 3) and diversity computation. The time complexity of the first two sub-routines are analyzed in Section C.1. In particular, the time cost of pattern enumeration is $O(2^{|E|})$. For the time cost of MNI computation, the node size of the extension pattern is $|V|$ in the worst case as the diversity measure does not has anti-monotonicity property. Thus, the time complexity to compute the MNI value of each candidate pattern is $O(|V|^{|V|})$, i.e., the cost of adapted GRAMI solution. For each candidate pattern, it costs $O(k \cdot |V|^2)$ to compute the diversity score among all these candidate patterns in top- k result set \mathcal{R} . Hence, the total time complexity of FPD-B is $O(2^{|E|} \cdot (|V|^{|V|} + k \cdot |V|^2))$.

Time Complexity of Algorithm FPD-A. Since the diversity score function does not have anti-monotonicity property, FPD-A will enumerate all possible candidate patterns, which is the same as FPD-B. Thus, the only difference of the time complexity between FPD-A and FPD-B is the join-based MNI computation. For each pattern, the time cost of join-based MNI computation is $O(|E|^{|E|})$ as the maximum number of edges of the candidate pattern is $|E|$. Hence, the total time complexity of kCP-A is $O(2^{|E|} \cdot (|E|^{|E|} + k \cdot |V|^2))$ with the diversity score computation cost analyzed above.

Space Complexity of Algorithm FPD-A. The space cost of FPD-A is similar with kCP-A but the maximum number of edges in candidate pattern is $|E|$, where $|E|$ is the cardinality of edge set in KG. Thus, the space complexity of FPD-A is $O(|E|^{|E|})$.

APPENDIX D

COMPLEMENTARY EXPERIMENT RESULTS

We report complementary experiment results for kCP and kCPD problems as follows.

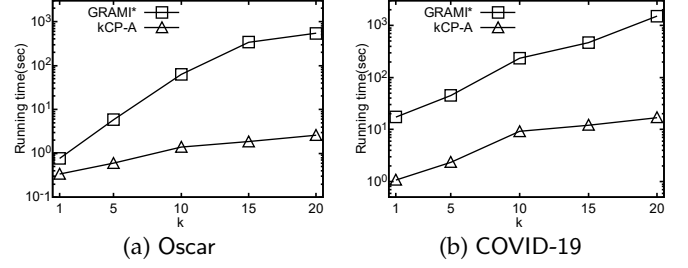


Fig. 1. Varying k : discover top- k frequent patterns without core pattern

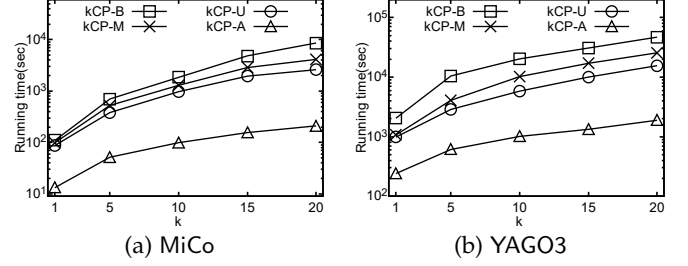


Fig. 2. Effect of k , kCP problem

D.1 Performance Evaluation on kCP Problem

Effect of k . (1) We adapted the state-of-the-art method GRAMI [3] to find top- k frequent patterns of the whole knowledge graph KG, denotes as GRAMI*. It is a baseline of our kCP problem with input core pattern $P = \emptyset$. For efficiency study, we revised kCP-A for finding top- k frequent patterns without a core, then compared it with GRAMI*. Figures 1(a) and (b) show the running time of kCP-A and GRAMI* on our default datasets by varying k . kCP-A achieves up to 210.2x and 89.5x faster than GRAMI* for Oscar and COVID-19, respectively. The experimental results confirm the superiority of kCP-A over its state-of-the-art method GRAMI*. (2) We present more efficiency results for large graphs, i.e., MiCo and YAGO3. Figure 2(a) reveals that the main performance gain comes from our join-based MNI computation instead of meta-index. In addition, the performance gap between kCP-A and its competitors (i.e., kCP-B, kCP-M and kCP-U) widens when k increases for all datasets, e.g., it achieves 24.8x on YAGO3, which suggests that kCP-A has good scalability on large datasets.

Memory Usage. We next present the memory usage of our algorithms. To measure the size of Java objects, we used a third-party tool *RamUsageEstimator* from Apache Lucene¹ in our implementation. Each reported memory is the peak memory that totally consumed by the algorithm. We calculate the memory usage (by counting the size of every constructed data structures) of the four solutions for kCP on the largest dataset YAGO3. As shown in Figure 3(a), the space consumptions of all methods are increasing with the number of frequent patterns k . kCP-A has the lowest memory usage among all methods as it employs the two-pass join scheme, which only materializes some of the pattern instances. In Figure 3 (b), we analyzed the space consumptions of different components in kCP-A by varying k . The intermediate table size of kCP-A increases with k as

1. <https://lucene.apache.org/>

kCP-A needs to compute MNI for more candidate patterns at large k .

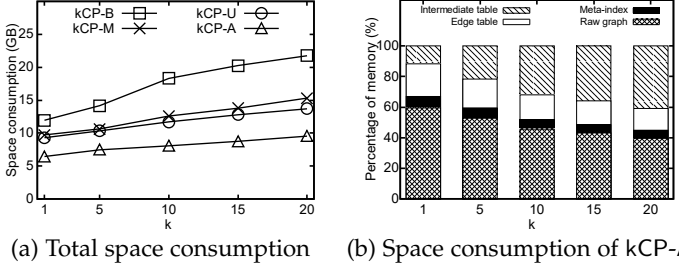


Fig. 3. Varying k : memory usage of kCP solutions, YAGO3

D.2 Performance Evaluation on kCPD Problem

Effect of k . We briefly present the experimental results of kCPD algorithms by varying k over large datasets. As shown in Figure 4, FPD-A performs the best among all the solutions, it achieves a speedup over one order of magnitude against FPD-B for MiCo and the largest dataset (i.e., YAGO3), which suggests that FPD-A has good scalability.

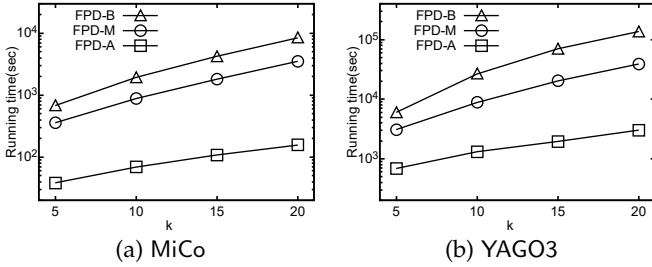


Fig. 4. Effect of k , kCPD problem

D.3 Case Study for kCPD Problem

We visualized the top-5 result patterns of GRAMI* (i.e., the top-5 most frequent patterns in KG), kCP (top-5 frequent patterns) and kCPD algorithms (top-5 frequent and diversified patterns) with core pattern P : Actor $\xrightarrow{\text{character is}}$ Character $\xrightarrow{\text{from}}$ Film, $\lambda = 0.3$ on the Oscar dataset in Figure 5. (I) the result patterns returned by GRAMI* are trivial. For example, 80% of the top-5 patterns returned by GRAMI* only includes one edge, however, the average edge size of the top-5 patterns produced by kCP and kCPD are 3 and 4. (II) kCPD discovers top-5 patterns with diversified node/edge types, e.g., its overlapping ratio (i.e., the average percentage of each node type exists in top- k patterns) is 12%, while the corresponding ratio on kCP and GRAMI* is 17% and 44%, respectively. (III) Compared to kCP problem, kCPD returns patterns that not only cover common entity types, e.g., "Actor" and "Award", but also involve uncommon entity types from the film crew, e.g., "Cinematographer". This provides more diversified options for Oscar knowledge graph exploration.

REFERENCES

- [1] M. Kuramochi and G. Karypis, "Finding frequent patterns in a large sparse graph," *DMKD*, vol. 11, no. 3, pp. 243–271, 2005.

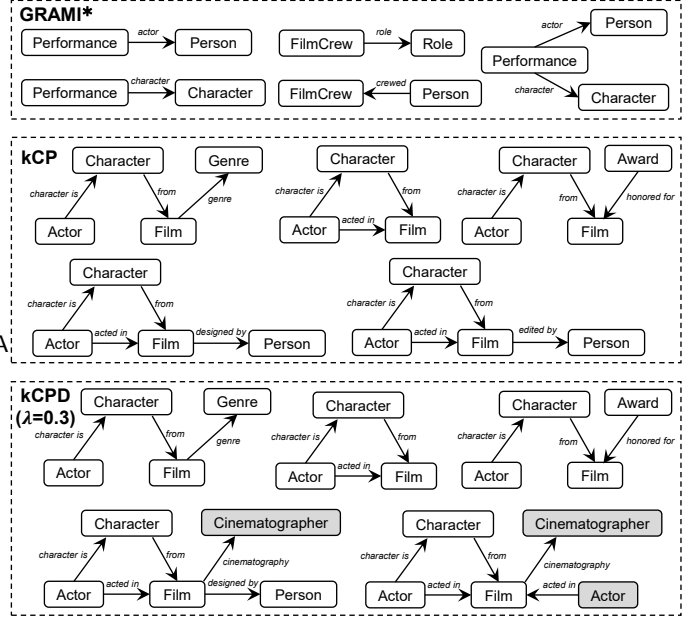


Fig. 5. Top-5 patterns returned by the three problems, Oscar

- [2] A. Prateek, A. Khan, A. Goyal, and S. Ranu, "Mining top-k pairs of correlated subgraphs in a large network," *PVLDB*, vol. 13, no. 9, pp. 1511–1524, 2020.
- [3] M. Elseidy, E. Abdelhamid, S. Skiadopoulos, and P. Kalnis, "Grami: Frequent subgraph and pattern mining in a single large graph," *PVLDB*, vol. 7, no. 7, pp. 517–528, 2014.